

مقایسه مدل‌های حاشیه‌ای برای تحلیل داده‌های طولی با پاسخ‌های دودویی

محتبی گنجعلی^۱، امید امدادی فر^۲

دانشگاه شهید بهشتی، دانشکده علوم ریاضی، گروه آمار

(دریافت: ۸۱/۴/۹؛ پذیرش: ۸۱/۸/۱۴)

چکیده

مدل حاشیه‌ای (Marginal model) یکی از رهابت‌هایی است که برای تحلیل داده‌های طولی (Longitudinal data) بکار می‌رود. در این مدل، همبستگی موجود بین پاسخ‌ها برای هر آزمودنی (Individual) به عنوان پارامتر درنظر گرفته می‌شود تا تحلیل‌های معنبری را نتیجه دهد. مدل‌های حاشیه‌ای مختلفی که برای تحلیل داده‌های طولی با پاسخ‌های دودویی ذکر شده‌اند از جمله مدل‌های حاشیه‌ای بالگ نسبت بخت حاشیه‌ای (Marginal odds ratio)، لک نسبت بخت شرطی (Conditional odds ratio)، مدل پروریت وابستگی (Dependence ratio)، مدل پروریت چند متغیره [Generalized Multivariate probit model] و روش معادلات برآوردگر تعمیم یافته (GEE) می‌رسند. مراور و مقایسه شده‌اند. مانده‌هایی برای بررسی یک‌کوئی برآشش این مدل‌ها معرفی گردیده‌اند. همچنین مدل‌های ذکر شده را در مثالی کاربردی برآشش داده‌ایم.

واژه‌های کلیدی: مطالعات طولی، مدل‌های حاشیه‌ای، معادلات برآوردگر تعمیم یافته، نسبت وابستگی، نسبت بخت حاشیه‌ای، نسبت بخت شرطی، مدل پروریت چند متغیره

۱- مقدمه

داده‌های طولی داده‌هایی هستند که در آنها پاسخ‌ها برای هر آزمودنی در طول زمان تکرار شده‌اند و چون اغلب بین پاسخ‌ها برای هر آزمودنی خاص همبستگی وجود دارد نیاز به روش‌های ویژه‌ای برای تحلیل وجود دارد تا بتوان استنباط‌های معتری به دست آورد. چندین رهیافت برای تحلیل چنین داده‌هایی بکار می‌روند که از جمله آنها می‌توان مدل‌های اثرات تصادفی (Transition models)، مدل‌های انتقالی (Random effects models) و مدل‌های حاشیه‌ای را

نام برد (Diggel, et al., 1994).

در بخش دوم مدل حاشیه‌ای را تعریف کرده، سپس انواع مدل‌های حاشیه‌ای را که برای داده‌های طولی دودویی بکار می‌روند در بخش‌های ۳ تا ۶ مرور کرده‌ایم. در بخش ۷ مقایسه‌ای کلی بین مدل‌های حاشیه‌ای ارائه شده، انجام داده‌ایم. از آنجا که در داده‌های طولی دودویی مانده‌هایی که بتوانند با درنظر گرفتن همبستگی پاسخ‌های یک آزمودنی خاص نیکوبی برآش مدل را انجام دهنده هنوز معرفی نشده است، مانده‌هایی برای مدل‌های ذکر شده در بخش ۸ معرفی کرده‌ایم. سپس مدل‌های بیان شده را در بخش ۹ بر روی داده‌های آسم روتنیتسکی و وپیج (Rotnitzky & Wypij, 1994) برآش داده‌ایم.

۲- مدل‌های حاشیه‌ای

مدل‌های حاشیه‌ای مدل‌هایی هستند که در آنها تاثیر متغیرهای تبیینی بر پاسخ‌ها به طور مجزا از همبستگی بین پاسخ‌ها برای یک آزمودنی معین مدل‌بندی می‌شود. در این نوع مدل‌ها، امید حاشیه‌ای $E(Y_{ij})$ به صورت تابعی از متغیرهای تبیینی مدل‌بندی می‌شود. امید حاشیه‌ای همان پارامتری است که در مطالعه مقطوعی نیز مدل‌بندی می‌شود (Agresti, 1990).

در یک مدل حاشیه‌ای امید حاشیه‌ای پاسخ $E(Y_{ij}) = \mu_{ij}$ به متغیرهای تبیینی x_{ij} به صورت $x'_{ij}\beta = h(\mu_{ij})$ وابسته است؛ که در آن h تابع پیوند (Link function) شناخته شده‌ای است. به عنوان مثال h برای پاسخ‌های دودویی، لوجیست می‌باشد. همچنین واریانس حاشیه‌ای به میانگین حاشیه‌ای به صورت $\text{var}(Y_{ij}) = v(\mu_{ij})\phi$ وابسته است که v تابع واریانس شناخته شده و ϕ پارامتر مقیاسی است که ممکن است، علاوه بر دیگر پارامترها نیاز به برآورد داشته باشد و همبستگی γ_{ij} و γ_{ik} به صورت تابعی از میانگین حاشیه‌ای آنها و پارامترهای اضافی α است یعنی:

$$\text{Corr}(Y_{ij}, Y_{ik}) = \rho(\mu_{ij}, \mu_{ik}, \alpha)$$

که (β) تابع شناخته شده‌ای از μ ‌ها است. ضرایب رگرسیونی حاشیه‌ای β – دارای تفسیری مشابه ضرایب رگرسیونی در تحلیل مقطعی هستند. می‌توان گفت مدل‌های حاشیه‌ای برای پاسخهای وابسته تعمیمی از مدل‌های خطی تعمیم یافته (Generalized linear models) برای پاسخهای مستقل می‌باشند.

۲- رهیافت بهادر (Bahadur)

بهادر (۱۹۶۱) توزیع توان (Y_{ij}) که در آن $Y_{ij} = (Y_{1j}, \dots, Y_{mj})$ برای $i = 1, \dots, m$ و $j = 1, \dots, n$ و پاسخهای دودویی هستند را به صورت زیر در نظر می‌گیرد.

$$\prod_{j=1}^n \mu_{ij}^{e_{ij}} (1 - \mu_{ij})^{1 - e_{ij}} \left(1 + \sum_{i < k} \rho_{ijk} e_{ij} e_{ik} + \sum_{j < k < l} \rho_{ijkl} e_{ij} e_{ik} e_{jl} e_{ml} + \dots + \rho_{i_1 i_2 \dots i_m} e_{i_1} e_{i_2} \dots e_{i_m} \right)$$

که در آن:

$$\rho_{i_1 i_2 \dots i_m} = E(e_{i_1} e_{i_2} \dots e_{i_m}), \quad \rho_{ijk} = E(e_{ij} e_{ik}), \quad e_{ij} = \frac{Y_{ij} - \mu_{ij}}{\{\mu_{ij} (1 - \mu_{ij})\}^{1/2}}$$

بنا براین توزیع توان پاسخ‌ها بوسیله فرضهایی درباره n -ها همبستگی‌های حاشیه‌ای

$$\rho_{ij} = (\rho_{i_1 i_2 \dots i_m}, \rho_{m-i_1 i_2 \dots i_m})'$$

مشخص می‌شود.

نمایش بهادر در بالا برای درجه‌های متفاوتی از وابستگی در میان Y ‌ها قابل بیان می‌باشد. مثلاً ساختار استقلال نتیجه می‌دهد که تمام همبستگی‌های حاشیه‌ای دودویی و مراتب بالاتر صفر است. اما همبستگی‌های حاشیه‌ای (که در مدل بهادر تعدادشان با افزایش n بسیار زیاد می‌شود) به صورت زیر بوسیله احتمال‌های حاشیه‌ای (میانگین‌های حاشیه‌ای که در پاسخ‌های دودویی احتمال‌های پیروزی حاشیه‌ای آند) محدود می‌شوند.

$$A \leq \text{Corr}(Y_{ij}, Y_{ik}) \leq B \quad (1)$$

که در آن:

$$A = \max \left\{ - \left(\frac{\mu_{ij} \mu_{ik}}{(1 - \mu_{ij})(1 - \mu_{ik})} \right)^{1/2}, \left(\frac{(1 - \mu_{ij})(1 - \mu_{ik})}{\mu_{ij} \mu_{ik}} \right)^{1/2} \right\}$$

$$B = \min \left\{ \left(\frac{\mu_{ij}(1-\mu_{ik})}{(1-\mu_{ij})\mu_{ik}} \right)^{\frac{1}{2}}, \left(\frac{(1-\mu_{ij})\mu_{ik}}{\mu_{ij}(1-\mu_{ik})} \right)^{\frac{1}{2}} \right\}$$

به همین دلیل بسیاری از محققین، معیارهایی برای همبستگی در نظر می‌گیرند که موجب می‌شود حدود همبستگی به میانگین‌ها محدود نباشد. آنها برای آزمودنی‌ام ($i = 1, \dots, m$) ($i = 1, \dots, m$) تابع چگالی توام پاسخ‌ها را به صورت زیر در نظر می‌گیرند:

$$f(y_i, \Psi_i, \Omega_i) = \exp \{ \Psi'_i y_i + \Omega'_i W_i - A(\Psi_i, \Omega_i) \} \quad (2)$$

که در آن:

$$W_i = (y_{i,1} y_{i,2}, \dots, y_{i,m-1} y_m, \dots, y_{i,1} y_{i,2} \dots y_m)$$

برداری $(1 \times n-1)$ از حاصلضرب‌های دوتایی و بالاتر $y_{i,1}, \dots, y_{i,m}$ و $\Psi_i = (\Psi_{i,1}, \dots, \Psi_{i,n})'$ و $\Omega_i = (\omega_{i,1}, \dots, \omega_{i,(n-1)m}, \dots, \omega_{i,12\dots n})'$ بردارهایی از پارامترهای متعارف می‌باشند. $A(\Psi_i, \Omega_i)$ مقدار ثابتی است که به y_i وابسته نیست.

۴- رهیافت‌های حاشیه‌ای دیگر براساس درستنمایی

در این بخش مدل‌های حاشیه‌ای براساس درستنمایی ارائه می‌شوند. در تمام این مدل‌ها، توزیع توام پاسخ‌ها بوسیله درنظر گرفتن فرض‌هایی به طور کامل مشخص می‌شوند. یکی از این فرض‌ها درباره ارتباط پاسخ‌ها است. در این بخش رهیافت‌های حاشیه‌ای که معیارهای همبستگی آنها محدود به احتمال‌های حاشیه‌ای نیستند – لگ نسبت بخت شرطی، نسبت وابستگی و لگ نسبت بخت حاشیه‌ای – بیان می‌شوند.

در تمام این مدل‌ها ارتباط پاسخ‌ها با متغیرهای کمکی بصورت $x_{ii} \beta = \log ii(\mu_{ii})$ و واریانس پاسخ‌ها بصورت $\text{var}(Y_{ii}) = \mu_{ii}(1-\mu_{ii})$ بیان می‌شود. تنها تفاوت آنها در معیار همبستگی پاسخ‌ها است. توجه داشته باشید که لوجیت ($\log ii$) به معنای لگاریتم بخت است یعنی:

$$\log ii(\mu_{ii}) = \log \left(\frac{P(Y_{ii} = 1)}{1 - P(Y_{ii} = 1)} \right)$$

۱-۴) مدل حاشیه‌ای با پارامترهای همبستگی به صورت توابعی از لگ نسبت بخت شرطی در این مدل‌بندی که توسط فیتز‌موریس و لرد پیشنهاد شده است (Fitzmaurice & Laird, 1993) در معادله ۲ همبستگی پاسخ‌ها برای هر آزمودنی براساس لگ نسبت بخت شرطی و توابعی از آنها قابل بیان است.

در معادله ۲ پارامترهای Ψ عباراتی از احتمال‌های شرطی به صورت:

$$\psi_y = \text{logit}\{\Pr(y_y = 1)\}$$

و پارامترهای Ω (تعریف شده در معادله ۲) عباراتی از لگ نسبت بخت شرطی است یعنی:

$$\exp(\omega_{rs}) = \frac{\Pr(Y_r = 1, Y_s = 1 | Y_u = 0, t \neq r, s)}{\Pr(Y_r = 1, Y_s = 0 | Y_u = 0, t \neq r, s)} \times \frac{\Pr(Y_r = 0, Y_s = 0 | Y_u = 0, t \neq r, s)}{\Pr(Y_r = 0, Y_s = 1 | Y_u = 0, t \neq r, s)}$$

و باید توجه کرد که μ تابعی از هر دو بردار Ψ و Ω است.

در این مدل پارامترهای همبستگی مقید به میانگین حاشیه‌ای نیستند و فضای پارامترها برای $\Omega_i = (\omega_{i1}, \omega_{i2}, \dots, \omega_{in})^T$ (فضای اقلیدسی $n-1$ -بعدی) می‌باشد. ولی این نوع پارامترها به n (تعداد پاسخ‌های تکرار شده برای هر آزمودنی) وابسته می‌باشند. یعنی می‌توان گفت این نوع مدل بندی برای داده‌هایی با تعداد پاسخ‌های تکرار شده متفاوت برای برخی از آزمودنی‌ها قابل کاربرد نیستند.

۴-۲) مدل حاشیه‌ای با پارامترهای همبستگی به صورت توابعی از نسبت وابستگی این مدل حاشیه‌ای توسط اخولم و دیگران، پیشنهاد شده است (Ekholm *et.al.*, 1995). در این روش استفاده از پارامتر میانگین (مولفه‌هایی که احتمال موفقیت توام تمام مراتب پاسخ‌ها را در بردارد) برای تحلیل رگرسیون پاسخ دودویی چند متغیره پیشنهاد شده است که در آن همبستگی را با استفاده از نسبتهای وابستگی (که عبارت‌هایی از پارامتر میانگین می‌باشد) بیان می‌کنند.

اگر $S(Y, W)$ را در نظر بگیریم که Y و W همان بردارهای تعریف شده در معادله ۲ هستند آنگاه به بردار مقادیر مورد انتظار $(S(Y, W))$ پارامتر میانگین ϕ می‌گویند (Barndorff-Nelson & Cox, 1994).

یعنی:

$$\phi = E\{S(Y)\} = (\phi_1, \dots, \phi_m, \phi_{12}, \dots, \phi_{(n-1)n}, \dots, \phi_{12\dots n})'$$

مولفه‌های ϕ احتمال‌های موفقیت توام حاشیه‌ای را برای اندیس‌های آن بیان می‌کند یعنی:

$$\phi_{11} = \Pr(Y_{11} = 1), \dots, \phi_{112} = \Pr(Y_{11} = 1, Y_{12} = 1), \dots,$$

$$\phi_{12\dots n} = \Pr(Y_{11} = 1, \dots, Y_{nn} = 1)$$

برای مثال پاسخ دودویی دومتغیره برای آزمودنی i ام، (Y_{i1}, Y_{i2}) ، با پارامتر میانگین ϕ_{i1}, ϕ_{i2} را در نظر بگیرید. نسبت وابستگی ω_{i12} ، و لگ نسبت وابستگی λ_{i12} ، به صورت زیر تعریف می‌شود:

$$\begin{aligned} \phi_{i12-n} &= \Pr(Y_{i1} = 1, \dots, Y_m = 1) \\ \omega_{i12} &= \frac{\phi_{i12}}{\phi_{i1}\phi_{i2}}, \quad \lambda_{i12} = \log(\omega_{i12}) \end{aligned} \quad (4)$$

که در آن $10 \times 10 - (\omega_{i12})$ میزان درصد بزرگی احتمال موفقیت توأم Y_{i1} و Y_{i2} در مقایسه با فرض استقلال را نشان می‌دهد.

نسبت‌های وابستگی مراتب بالاتر نیز به صورت زیر تعریف می‌شود:

$$\omega_{i123} = \frac{\phi_{i123}}{\phi_{i1}\phi_{i2}\phi_{i3}}, \dots, \omega_{i123-n} = \frac{\phi_{i123-n}}{\phi_i\phi_{i2}\dots\phi_m}.$$

نسبت‌های وابستگی بوسیله مقایسه احتمال‌های موفقیت حاشیه‌ای با حالت استقلال تفسیر می‌شوند، اگر نسبت وابستگی با دو اندیس و بیشتر، برابر با یک باشد- برای نمونه اگر $\omega_{i123} = \phi_{i1}\phi_{i2}\phi_{i3} = 1$ آنگاه $\phi_{i123} = \phi_{i1}\phi_{i2}\phi_{i3}$ که استقلال سه‌تایی (Y_{i1}, Y_{i2}, Y_{i3}) را نتیجه می‌دهد. باید توجه داشت که اگر موفقیت و شکست برای مولفه‌هایی دو دوی معاوضه شوند (مثلاً برای (Y_{i2}, Y_{i1}) و (Y_{i1}, Y_{i2})) مطالعه شود آنگاه ضریب همبستگی و لگ نسبت بخت در علامت تغییر می‌کنند اما در مقدار قدر مطلق تغییر نمی‌کنند اما ω_{i123} به $\frac{1 - \phi_{i12}\omega_{i12}}{1 - \phi_{i12}}$ تبدیل می‌شود. علاوه بر این ضریب همبستگی و نسبت‌های بخت برای (Y_{i1}, Y_{i2}) و (Y_{i2}, Y_{i1}) برابر با (Y_{i1}, Y_{i2}) است در حالی که نسبت وابستگی برای (Y_{i2}, Y_{i1}) و (Y_{i1}, Y_{i2}) تنها تابعی صعودی از نسبت وابستگی برای (Y_{i1}, Y_{i2}) است.

۳-۴) مدل حاشیه‌ای با لگ نسبت بخت حاشیه‌ای به عنوان پارامتر همبستگی: این مدل بندهی توسط لیانگ و دیگران و لیپسیتز و دیگران ارائه شده است (Liang et al., 1992; Lipsitz et al., 1991). در این مدل بندهی، در معادله ۳ پارامترهای Ψ عباراتی از احتمال‌های حاشیه‌ای می‌باشند یعنی:

$$\psi_{\eta} = \log it\{\Pr(Y_{\eta} = 1)\}$$

و پارامترهای دو طرفه، Ω لگ نسبت بخت حاشیه‌ای می‌باشند یعنی:

$$\exp(\omega_{ns}) = OR(Y_{nr}, Y_{ns}) = \frac{\Pr(Y_{nr} = 1, Y_{ns} = 1)}{\Pr(Y_{nr} = 0, Y_{ns} = 1)} \frac{\Pr(Y_{nr} = 0, Y_{ns} = 0)}{\Pr(Y_{nr} = 1, Y_{ns} = 0)}$$

و پارامترهای سه طرفه و بالاتر به صورت توابعی از لگ نسبت بخت حاشیه‌ای است. به عنوان مثال:

$$\exp(\omega_{rst}) = \frac{OR(Y_{nr}, Y_{ns} | Y_{rl} = 1)}{OR(Y_{nr}, Y_{ns} | Y_{rl} = 0)}$$

یا بطور معادل:

$$\omega_{rst} = \log[OR(Y_{nr}, Y_{ns} | Y_{rl} = 1)] - \log[OR(Y_{nr}, Y_{ns} | Y_{rl} = 0)]$$

است.

با این نوع مدل‌بندی‌ها، درجه‌های مختلف واپستگی در میان Y ‌ها قابل بیان است. به عنوان مثال مدل استفلال نتیجه می‌دهد که تمام پارامترهای همبستگی دو طرفه و بالاتر صفر است و یا اگر تمام عناصر Ω مخالف صفر باشند مدل انتساب شده می‌باشد. بین این دو حالت مدل‌های مختلفی از واپستگی را می‌توانیم در نظر بگیریم. همانطور که ملاحظه می‌شود اینگونه مدل‌ها محاسبات زیادی برای برآورد پارامترها دارند (جز برای $n=2$). در نتیجه برای رهایی از این مشکل محققین اکثراً ارتباطهای مراتب بالاتر از دو را صفر در نظر می‌گیرند (Zhao & Prentice, 1990).

۵- مدل‌های حاشیه‌ای با استفاده از مفهوم متغیر پنهان

از مدل پربویت چند متغیره با استفاده از مفهوم متغیر پنهان (Latent variable) برای تحلیل‌های طولی دو دویی استفاده می‌شود (Long, 1997) چون Y^* متغیری پیوسته است دسته‌ای از مدل‌های خطی را می‌توان برای تحلیل آنها در نظر گرفت، که به صورت زیر تعریف می‌شوند.

$$Y^*_{nl} = X'_{nl}\beta + \varepsilon_{nl}$$

که در آن فرض‌هایی درباره ε ‌ها در نظر گرفته می‌شود (Long, 1997). به عنوان مثال: $\varepsilon_i = (\varepsilon_{i1}, \dots, \varepsilon_{in})'$ برای آزمودنی‌ها، دارای توزیع نرمال n متغیره با میانگین صفر و

$$\Sigma = \begin{pmatrix} 1 & \rho_{12} & \dots & \rho_{1n} \\ \rho_{21} & 1 & \dots & \rho_{2n} \\ \vdots & & \ddots & \\ \rho_{n1} & \dots & \dots & 1 \end{pmatrix} = \rho$$

درنظر گرفته می‌شود. این فرض مدل پروریت چند متغیره را نتیجه می‌دهد. با توجه به این که ابعاد مشاهده شده بر اساس قاعده زیر مشخص می‌شوند.

$$y_n = \begin{cases} 1 & y_n^* > \tau \\ 0 & y_n^* \leq \tau \end{cases}$$

که $\tau = 0$ اختیار می‌شود (Long, 1997)، احتمال‌های توام و درنتیجهتابع درستنمایی را با استفاده از فرض‌هایی روی ϵ می‌توان محاسبه نمود.

همانطور که ملاحظه می‌شود با افزایش n تعداد پارامترهای ماتریس همبستگی افزایش چشمگیری خواهد داشت، $\frac{n(n-1)}{2}$ پارامتر، که برآورد تمام این پارامترها نیاز به محاسبات وقت زیاد دارد. به همین دلیل برای ماتریس همبستگی، حالت‌های خاصی درنظر گرفته می‌شود. به عنوان مثال:

$$Corr(\varepsilon_{it}, \varepsilon_{it'}) = \rho, \forall t \neq t' \quad (1)$$

یا

$$Corr(\varepsilon_{it}, \varepsilon_{it'}) = \rho^{|t-t'|} \quad (2)$$

که در حالت ۲ همبستگی بین دو پاسخ آزمودنی یکسان با افزایش فاصله زمانی بین دو پاسخ کاهش می‌یابد.

۶- رهیافت معادلات برآورد تعمیم یافته

این رهیافت توسط لیانگ و زیگر پیشنهاد شده است (Liang and Zeger, 1981). رهیافت معادلات برآورد تعمیم یافته، رهیافت واحدی را برای تحلیل برآمدهای گوناگون پیوسته و گستته پیشنهاد می‌کند. در این میان دسته‌ای از معادلات برآوردگر تعمیم یافته (GEE) برای پارامترهای رگرسیون پیشنهاد شده است. در این رهیافت تنها فرض‌هایی که اختیار می‌شود فرض درباره امید حاشیه‌ای مرتبه اول و دوم پاسخ‌ها است و درباره توزیع کامل آنها فرضی اختیار نمی‌شود.

این رهیافت برای پارامترهای رگرسیون برآوردهای سازگار می‌دهد. دراصل (GEE) حالت چند متغیره شبه درستنمایی (Quasi likelihood) است که توسط ودربرن (Wedderburn, 1974) بیان شده است.

در این رهیافت تابعی از امید حاسیه‌ای متغیر وابسته به صورت تابعی خطی از متغیرهای کمکی مشخص می‌شود. همچنین فرض می‌شود که واریانس تابعی معلوم از میانگین حاسیه‌ای است. به علاوه ماتریس همبستگی عملی برای مشاهدات هر آزمودنی، معین می‌شود. یعنی:

$$\text{var}(Y_{ii}) = v(\mu_{ii}) \phi / h(\mu_{ii}) = x'_{ii} \beta$$

و همبستگی عملی مشاهدات برای آزمودنی i ام (α_i) فرض می‌شود که به بردار s از پارامترهای اضافی، α ، وابسته است. α ، همبستگی پاسخ‌های آزمودنی i ام را به طور کامل بیان می‌کند. اکثراً از (α) R به عنوان ماتریس همبستگی عملی استفاده می‌شود، زیرا انتظار این که ماتریس همبستگی واقعی به درستی مشخص شود نیست. درنتیجه ماتریس کوواریانس عملی V_i به صورت:

$$V_i = A_i^{\frac{1}{2}} R_i (\alpha) A_i^{\frac{1}{2}} / \phi$$

است که A_i ماتریس قطری، n_i با عناصر (μ_{ii}) در زامین عنصر قطری آن می‌باشد.
معادلات برآورد به صورت زیر می‌باشد.

$$\sum_{i=1}^m D_i' V_i^{-1} S_i = 0 \quad (4)$$

که در آن $D_i = \frac{\partial \mu_i}{\partial \beta}$ و $S_i = Y_i - \mu_i$ است. هنگامی که $n_i = 1, \dots, m$ معادله (4)

به معادله شبه درستنمایی تحویل می‌شود (Wedderburn, 1974). معادله ۴ می‌تواند به صورت تنها تابعی از β بیان شود به این صورت که برآوردهای سازگار α و ϕ در معادله ۴ قرار گیرد (لینگ و زیگر، ۱۹۸۶) یعنی به صورت زیر خواهد بود.

$$\sum_{i=1}^m U_i [\beta, \dot{\alpha} \{ \beta, \dot{\phi}(\beta) \}] = 0 \quad (5)$$

$\hat{\beta}_i$ که نشان دهنده برآوردهای β تحت رهیافت (GEE) است، جواب معادله بالا خواهد بود. براساس قضیه ۲ در لینگ و زیگر (۱۹۸۶)، $m^{\frac{1}{2}} (\hat{\beta}_i - \beta)$ بطور مجانبی گاووسی با میانگین صفر و ماتریس کوواریانس V_i به صورت زیر است:

$$\begin{aligned} V_{G_i} &= \lim_{m \rightarrow \infty} m \left(\sum_{i=1}^m D_i' V_i^{-1} D_i \right)^{-1} \\ &\times \left\{ \sum_{i=1}^m D_i' V_i^{-1} \text{cov}(Y_i) V_i^{-1} D_i \right\} \left(\sum_{i=1}^m D_i' V_i^{-1} D_i \right)^{-1} \\ &= \lim_{m \rightarrow \infty} m V_i^{-1} \hat{V}_{G_i} \end{aligned}$$

برآورد واریانس $\hat{\beta}_i$ می‌تواند برای یافتن واریانس برآورده‌گرها استفاده شود
(Zeger & Liang, 1992)

۷- مقایسه کلی مدل‌های حاشیه‌ای

در این بخش به بررسی و مقایسه مدل‌های حاشیه‌ای ارائه شده دربخش‌های قبل می‌پردازیم. ابتدا مدل استقلال، که فرض بر عدم همبستگی بین پاسخ‌های متعلق یک آزمودنی خاص دارد، را در نظر گیرید. براساس قضیه ۱ لیانگ و زیگر (۱۹۸۶)، $\hat{\beta}_i$ (برآورد پارامتر تحت مدل استقلال) برآورده سازگار برای β مهیا می‌کند. اما این روش استنباط‌هایی نادرست درباره پارامترها می‌دهد زیرا فیتزموریس و دیگران (۱۹۹۳) بیان کردند که اکثرا خطای استاندارد برآورده پارامترهای مربوط به متغیرهای زمان مانع میل به کم برآورد شدن و برای پارامترهای مربوط به متغیرهای زمان متغیر میل به بیش برآورد شدن دارند. و بنابراین برآوردهای واریانس‌ها سازگار نیستند.

رهیافت‌هایی که دربخش چهارم بیان شدند هرکدام دارای محسن و معایبی می‌باشند. پارامترهای همبستگی به صورت توابعی از لگ نسبت بخت شرطی بر خلاف همبستگی‌های بهادر و لگ نسبت بخت حاشیه‌ای مقید به میانگین حاشیه‌ای نیستند. ولی این نوع پارامترها به n (تعداد پاسخ‌های تکرار شده برای هر آزمودنی) وابسته می‌باشند.

وقتی پارامترهای همبستگی به صورت توابعی از لگ نسبت بخت حاشیه‌ای بیان شود، پارامترهای همبستگی به n وابسته نیست. و بهمین دلیل برای تعداد پاسخ‌های تکرار شده متفاوت برای آزمودنی‌های مختلف، این مدل کاربرد دارد. دیگر مزیت مدل لگ نسبت بخت حاشیه‌ای نسبت به مدل با لگ نسبت بخت شرطی آن است که قابل دوباره تولید شدن می‌باشد ولی مدل با لگ نسبت بخت شرطی قابلیت دوباره تولید شدن را ندارد (مدلی قابل دوباره تولید شدن است که اگر $(Y_{i1}, Y_{i2}, \dots, Y_{im})'$ دریک مدل حاشیه‌ای صدق کند، آنگاه زیر مجموعه‌ای

از Y مانند $(Y_{i1}', Y_{i2}', \dots, Y_{im}')$ نیز در مدل با پارامترهای متناظر آن صدق کند).

هنگامی که پارامترهای همبستگی بین پاسخ‌ها براساس نسبت وابستگی بیان می‌شود، پارامترهای همبستگی همانند لگ نسبت بخت حاشیه‌ای وابسته به β نیست. باید توجه داشت که پیچیدگی قید پارامترهای همبستگی نسبت‌های وابستگی ولگ نسبت بخت حاشیه‌ای به میانگین‌ها کمتر از پیچیدگی قید همبستگی‌های بهادر است.

در هر سه مدل بخش^۴، ضرایب رگرسیونی و پارامترهای همبستگی بر هم عمود می‌باشند. این بدهی معنی است که کواریانس برآوردهای ضرایب رگرسیونی و پارامترهای همبستگی برابر صفر است. (بارندورف نیلسون و کاکس، ۱۹۹۴).

برآوردهای ماکسیمم درستنمایی پارامترهای مدل‌های ذکر شده در بخش چهارم براساس روش‌های تکراری عددی بدست می‌آیند. برای تمام این مدل‌ها بدون توجه به ساختار همبستگی، برآورد پارامترهای رگرسیونی سازگار می‌باشد ولی تنها اگر ساختار همبستگی برای این مدل‌ها بدرستی مشخص شده باشد آنگاه برآورد (β) نیز برآوردهای سازگار می‌باشد.

رهیافت GEE که رهیافتی غیردرستنمایی است برخلاف مدل پربویت چند متغیره که به فرض توزیع چند متغیره نرمال برای متغیرهای پنهان نیاز دارد، تنها به فرضهایی درباره دو گشتاور اول پاسخ‌ها نیاز دارد. البته باید توجه کرد که در توزیع گاوی، فرض درباره دو گشتاور اول، توزیع را به طور کامل مشخص می‌کند در حالی که برای داده‌های دودویی که مورد بحث اصلی است، این مطلب برقرار نمی‌باشد (توجه کنید که متغیرهای پنهان ما دارای توزیع گوسی‌اند نه متغیرهای مشاهده شده). در GEEتابع پیوندهای مختلفی همچون پربویت، لوجیت و ... می‌تواند مورد استفاده قرار گیرد، ولی در پربویت چند متغیره تنها پیوند پربویت و در مدل برای لگ نسبت بخت حاشیه‌ای تنها پیوند لوجیت مورد استفاده است. GEE همانند مدل با نسبت وابستگی و نسبت بخت حاشیه‌ای، برای داده‌های با اندازه تکرارهای متفاوت برای آزمودنی‌های مختلف قابل اجرا است.

یکی از محدودیتهای GEE نداشتن روش‌هایی برای آزمون نیکویی برآذش، مقایسه مدل و استنباط بر روی پارامترها بر اساس درستنمایی است زیرا در این روش، توزیع توام به طور کامل مشخص نمی‌شود. با توجه به این که برآوردهای دارای توزیع نرمال مجانی هستند، در این روش برای حل این مشکل از آماره‌های والد (Wald Statistics) استفاده می‌شود. یکی دیگر از محدودیتهای GEE آن است که اگر مدل برای میانگین‌های حاشیه‌ای به درستی مشخص نشده باشد، برآوردهای ناسازگار برای β به دست می‌آید.

۸- بررسی مانده‌ها

بررسی مانده‌ها برای نیکوبی برازش مدل و یافتن دورافتاده‌ها (مشاهداتی که روی نیکوبی برآذش مدل تاثیر می‌گذارند) در حالت مطالعات طولی، همانند مطالعات مقطعی ضروری است. تفاوتی که در این حالت برای مانده‌ها نسبت به مطالعات مقطعی وجود دارد این است که برای هر آزمودنی چند پاسخ داریم. برای محاسبه مانده‌ها اگر $\{\beta\} = \text{diag}\{\text{var}(Y_i)\}$, $\Sigma = \text{DrN}$ گرفته شود آنگاه همبستگی بین پاسخ‌های متعلق به یک آزمودنی خاص درنظر گرفته نمی‌شود. اما اگر از ماتریس کوواریانس Y استفاده شود، این همبستگی برای محاسبه مانده‌ها نیز مورد استفاده قرار می‌گیرد. در این حالت مانده‌های پیرسون برای مشاهدات Y را به صورت زیر معرفی می‌کنیم. مانده‌ها به صورت:

$$r_i^p = \sum_i^{-1/2} (\hat{\beta})(y_i - \hat{\mu}_i)$$

داده می‌شود که در آن $(\hat{\beta})$ برآورد ماتریس کوواریانس مشاهدات است. همچنین آماره نیکوبی برازش پیرسون به صورت:

$$\chi^2 = \sum_{i=1}^m \chi_p^2(y_i, \hat{\mu}_i)$$

است که در آن:

$$\chi_p^2(y_i, \hat{\mu}_i) = (y_i - \hat{\mu}_i)' \Sigma^{-1} (\hat{\beta})(y_i - \hat{\mu}_i)$$

$\chi_p^2(y_i, \hat{\mu}_i)$ را می‌توان به صورت: $r_i^{p'} r_i^{p'}$ که در آن $r_i^{p'}$ تراهنده بردار r_i^p تعریف شده در بالا است، نوشت. بزرگ بودن χ_p^2 در مقایسه با درجه‌آزادی آن (تعداد پارامترهای برآورد شده- m) گواه بر بدی برآذش مدل دارد. در این صورت یا مدل مورد استفاده مناسب نیست یا، در صورت وجود، مؤلفه‌های دیگری بایستی در مدل به عنوان متغیرهای تبیینی درنظر گرفته شود. به عنوان مثال با مدل دو متغیره پروبیت با دو پاسخ برای هر آزمودنی به جای ساختار کوواریانس زیر که فرض استقلال را درنظر می‌گیرد:

$$\Sigma(\beta) = \begin{bmatrix} \pi_{11}(1-\pi_{11}) & \cdot \\ \cdot & \pi_{12}(1-\pi_{12}) \end{bmatrix}$$

ماتریس کوواریانس Y به صورت زیر را در نظر می‌گیریم:

$$\Sigma(\beta) = \begin{bmatrix} \pi_{11}(1-\pi_{11}) & \pi_{12}-\pi_{11}\pi_{12} \\ \pi_{12}-\pi_{11}\pi_{12} & \pi_{12}(1-\pi_{12}) \end{bmatrix}$$

که در آن:
www.SID.ir

$$\text{cov}(Y_{i1}, Y_{i2}) = \pi_{i12} - \pi_{i1}\pi_{i2}$$

$$\pi_{i12} = \Pr(Y_{i1} = 1, Y_{i2} = 1)$$

$$= 1 - \Phi(-x'_{i1}\beta_1) - \Phi(-x'_{i2}\beta_2) + \Phi_2(-x'_{i1}\beta_1, -x'_{i2}\beta_2, \rho)$$

$$\pi_{ij} = \Pr(Y_{ij} = 1) = \Phi(x'_{ij}\beta) \quad j=1, 2.$$

با

و

(۲۰) تابع توزیع نرمال دو متغیره با همبستگی ρ بین دو متغیر است $\Phi(\cdot, \cdot, \rho)$ تابع توزیع نرمال استاندارد است. اگر برآورد β را در $\sum_i (\beta)^{-1/2}$ قرار دهیم برآورد این ماتریس بدست می‌آید. همان طور که ملاحظه می‌شود برآورد پارامتر ρ در محاسبه برآورد ماتریس کوواریانس و بنابراین مانده‌ها نقش دارد.

۹- مثال عملی: (برازش مدل‌های حاشیه‌ای بر روی داده‌های آسم)

این داده‌ها (رتیتزرکی و واپیچ، ۱۹۹۴) در ۶ شهرستان منطقه هاروارد، ایالو جمع‌آوری شده‌اند، به این ترتیب که ۷۰۶ پسر و ۷۱۳ دختر سفید پوست را در ۹ سالگی و دوباره در ۱۳ سالگی از نظر ابتلا یا عدم ابتلا به آسم مورد مطالعه قرار داده‌اند. می‌خواهیم تاثیر زمان و جنس را بر احتمال داشتن آسم مورد آزمون قرار دهیم. جدول ۱ داده‌های کامل و بدون مقادیر گمشده برای متغیر پاسخ در این داده‌ها را نشان می‌دهد. در بررسی‌های جدا برای رد آزمون کاملاً تصادفی بودن مقادیر گمشده گواهی بدست نیامد (Little and Rubin, 1987). همانطور که جدول ۱ نشان می‌دهد ۵۵۷ نفر پسر و ۵۹۰ نفر دختر به متغیر مورد علاقه پاسخ داده‌اند. برای بررسی تاثیر سن و زمان بر احتمال داشتن آسم، مدل‌های حاشیه‌ای که قبل ذکر شد را برازش می‌دهیم.

مدل برای نسبت وابستگی، لگ نسبت بخت حاشیه‌ای و مدل استقلال به صورت زیر در نظر گرفته شده است:

$$\log it(\pi_{ij}) = a_i + b_j sex, \quad j=1, 2, i=1, \dots, 1147$$

که در آن π_{ij} احتمال داشتن آسم فرد i ام در زمان j ام است، sex جنس آزمودنی i ام (1 : پسر، 0 : دختر) را نشان می‌دهد و $\text{var}(Y_{ij}) = \pi_{ij}(1 - \pi_{ij})$ است. پارامتر همبستگی برای نسبت وابستگی با لگ نسبت وابستگی و با لگ نسبت بخت حاشیه‌ای برای

مدل نسبت بخت حاشیه‌ای بیان می‌شود. مدل حاشیه‌ای پربویت دو متغیره با مفهوم متغیر پنهان نیز به این داده‌ها برآذش داده‌ایم. مدل در نظر گرفته شده به صورت:

$$Y_{ij}^* = a_j + b_j \text{sex}_i + \varepsilon_{ij} \quad j=1,2$$

است. که Y_{1j}^* و Y_{2j}^* متغیرهای پنهان هستند. همچنین همانطور که درخشش ۵ ذکر شد $\varepsilon_{ij} = (\varepsilon_{1j}, \varepsilon_{2j})'$ برداری با توزیع نرمال دو متغیره است.

جدول ۱ - داده‌های مربوط به آسم

		سالگی ۹			
		۱۳ سالگی			
پسرها	دخترها	اسم دارد	اسم ندارد	جمع کل	
		۲۲	۶	۲۸	
		۱۵	۵۱۴	۵۲۹	
		جمع کل	۳۷	۵۲۰	
				۵۵۷	
		اسم دارد	۱۳	۳	
		اسم ندارد	۱۳	۵۶۱	
		جمع کل	۲۶	۵۶۴	
				۵۹۰	

نتایج برآذش مدل‌ها در جدول ۲ آمده است. با استفاده از نتایج بدست آمده در جدول ۲ مشاهده می‌شود که پارامترهای همبستگی وجود همبستگی بالای مثبت ($\rho = 0.000$) را بین دو پاسخ برای هر آزمودنی نشان می‌دهد.

همان‌گونه که ملاحظه می‌شود تمام این مدل‌ها ضریب ثابت را معنی‌دار نشان می‌دهند یعنی عواملی مانند زمان بوده‌اند که درنظر نگرفته‌ایم. به این دلیل دوباره مدل GEE را بر روی داده‌ها برآذش داده‌ایم. مدل GEE به صورت زیر تعیین می‌گردد:

$$\log it(\pi_{ij}) = a_1 + b_1 \text{sex}_i + cI(age13)_{ij} \quad j=1,2$$

که در آن $(age13)_1$ به صورت زیر تعریف می‌شود:

$$I(age13) = \begin{cases} 1 & \text{اگر فرد سیزده ساله باشد،} \\ & \text{اگر فرد سیزده ساله نباشد،} \\ 0 & \end{cases}$$

در GEE ماتریس همبستگی به صورت غیر ساختاری درنظر گرفته شده است. نتایج برآزش مدل‌ها در جدول ۲ آمده است.

هر چند مقایسه مدل‌ها در یک مثال امکان پذیر نیست و نیاز به بررسی‌های نظری الزامی است، خلاصه‌ای از عملکرد مدل‌های مختلف در برآورد پارامترها می‌آوریم. در همه مدل‌ها گواه کافی بر معنی داری پارامتر منتبه به زمان (c یا تفاضل $a_1 - a_2$) وجود دارد. به این معنی که هر چه زمان اقامت در این ۶ شهرستان بیشتر می‌شود احتمال ابتلاء به آسم افزایش می‌باید. در مدل حاسیه‌ای پربویت دو متغیره جنس در زمان اول معنی دار است ولی در زمان دوم معنی دار نیست. به این معنی که پسران در سنین پایین برای آسم داشتن محتمل‌تر از دختران هستند ولی با گذشت زمان پسران بهبود یافته بیشتر از دخترانی هستند که به آسم مبتلا می‌شوند. مدل پربویت با فرض استقلال نیز به نتیجه یکسانی می‌رسد، اما واریانس‌های برآوردگرهای پارامترها کمی بیش یا کم برآورد می‌شوند. مدل استقلال با پیوند لوچیت نیز در مورد برآورد پارامترها به نتیجه‌ای شبیه به مدل پربویت دو متغیره می‌رسد، ولی این مدل در مقایسه با مدل نسبت واپسگی با پیوند لوچیت یا نسبت بخت با پیوند لوچیت واریانس برآوردگرهای را بسیار متزلزل برآورد می‌کند.

جدول ۲- برآورد پارامترها تحت مدل‌های در نظر گرفته شده.

پارامتر	نسبت بخت با پیوند لوچیت		نسبت واپسگی با پیوند لوچیت		مدل استقلال با پیوند لوچیت	
	خطای استاندارد		خطای استاندارد		خطای استاندارد	
	خطای استاندارد	برآورد	خطای استاندارد	برآورد	خطای استاندارد	برآورد
a_1	-۰/۵۰	۰/۴۱۸	-۳/۷۹۰	۰/۲۳۲	-۳/۱۸۰	۰/۱۲۳
a_2	-۳/۷۲۳	۰/۰۲۶	-۲/۸۸۸	۰/۱۷۱	-۳/۰۷۷	۰/۱۲۰
b_1	۰/۰۵۲۴	۰/۰۴۱۳	۰/۰۲۷۸	۰/۰۲۶۲	۰/۰۶۴۱	۰/۰۳۱۹
b_2	۰/۰۱۵۳	۰/۰۳۴۲	۰/۰۰۷۳	۰/۰۱۸۶	۰/۰۴۳۴	۰/۰۲۶۳
پارامتر همبستگی	۴/۹۸۶	۰/۰۴۲۱	۱۴/۰۴۱	۱/۶۴۵	----	----
منهای لگاریتم درستنمایی	۳۳۷/۲۱۱		۳۳۸/۷۴۶		۴۲۷/۱۸۱	
پارامتر	مدل GEE با پیوند لوچیت		مدل پربویت		مدل استقلال با پیوند پربویت	
	خطای استاندارد		خطای استاندارد		خطای استاندارد	
	خطای استاندارد	برآورد	خطای استاندارد	برآورد	خطای استاندارد	برآورد
a_1	-۳/۴۸۶	۰/۲۱۱	-۱/۹۲۵	۰/۱۰۵	-۱/۹۲۷	۰/۱۰۶
a_2	----	----	-۱/۷۰۵	۰/۰۹۰	-۱/۷۰۵	۰/۰۹۱
b_1	۰/۰۴۸۹	۰/۰۲۵۵	۰/۰۲۸۳	۰/۰۱۳۵	۰/۰۲۸۳	۰/۰۱۳۹
b_2	----	----	۰/۰۲۰	۰/۰۱۲۱	۰/۰۲۰۲	۰/۰۱۲۲
c	۰/۰۳۷۷	۰/۰۱۲۱	----	----	----	----
پارامتر همبستگی	۰/۰۶۴۴	----	۰/۰۹۲۶	۰/۰۲۴۵	----	----
منهای لگاریتم درستنمایی	----		۳۳۷/۱۰۶		۴۲۷/۱۸۱	

در مدل‌های نسبت وابستگی با پیوند لوجیت یا نسبت بخت با پیوند لوجیت اثر جنس در هیچیک از زمانها معنی‌دار نیست، به این معنی که پسرها و دخترها در هر سنی به یک‌اندازه برای داشتن آسم محتملند. در مدل GEE، که در آن بر خلاف مدل‌های دیگر هیچ فرضی در مورد توزیع پاسخها اعمال نمی‌شود، گواهی قوی بر معنی‌داری اثر زمان ولی گواهی ضعیف بر معنی‌داری جنس وجود دارد.

- برای نیکوبی برازش مدل بر اساس مانده‌ها که در بخش ۸ بیان شده است، می‌توانیم عمل کنیم. مثلاً برای مدل پربویت دو متغیره برازش شده بر روی این داده‌ها، آماره کی دو برابر با $2542/0.89$ است که با این مقدار آماره فرض مناسب بودن مدل رد می‌شود. اگر در مدل پربویت دو متغیره همبستگی را در نظر نگیریم مقدار آماره برابر با $2558/0.84$ است و خوبی برازش این مدل نیز رد می‌شود. فکر می‌کنیم که عدم برازش خوب به دلیل وجود عامل‌های متعدد موثر دیگری همچون میزان آسودگی هوا و عوامل محیطی و ژنتیکی دیگر است، که در این بررسی در نظر گرفته نشده است. بهر حال کاهش مقدار آماره کی دو در مدل عدم استقلال نسبت به مدل استقلال به خاطر در نظر گرفتن همبستگی بین پاسخ‌ها است.

Reference

- Agresti, A. (2002) *Categorical Data Analysis*. New York: John Wiley.
- Bahadur, R.R. (1961) *A representation of the joint distribution of responses to n dichotomous item. In studies on item analysis and prediction*. Stanford: Standorf University Press.
- Barndorff-Nielson, O.E., and Cox, D.R. (1994) *Inference and Asymptotics*. London: Chapman & Hall.
- Diggle, P.J., Liang, K., Zeger, S. (1994) *Analysis of Longitudinal Data*. Oxford Science Publication.
- Ekholm, A., Smith, P.W.F., McDonald, J.W. (1995) Marginal regression analysis of a multivariate binary response. *Biometrika*, **82**, 847-854.
- Fitzmaurice, G.M., Laird, N.M. (1993) *A Likelihood-Based method for analysing longitudinal binary responses*. *Biometrika*, **80**, 141-151.
- Fitzmaurice, G.M., Laird, N.M. Rotnitzky, A.G. (1993) *Regression Models for Discrete Longitudinal responses*. *Statistical Science*, **8**, 284-309.

- Liang, K.Y., Zeger, S.L. (1986) *Longitudinal data analysis using generalized linear models*. Biometrika, **73**, 13-22.
- Liang, K.Y., Zeger, S.L., Qaqish, B. (1992) *Multivariate regression analysis for categorical data (with discussion)*. Journal of the Royal Statistical Society, B, **54**, 3-40.
- Lipsitze, S., Laird, N., Harrington, D. (1991) *Generalized estimating equations for correlated binary data: Using odds ratios as a measure of association*. Biometrika, **78**, 153-160.
- Little, R. J. and Rubin, D. (1987) *Statistical analysis with missing data*. New York: Wiely.
- Long, S.J. (1997) *Regression models for categorical and limited dependent variables*. London: SAGE.
- Rotnitzky, A. and Wypij, D. (1994) *A note on the Bias of Estimation with Missing Data*. Biometrika, **147**, 87-99.
- Wedderburn, R.W.M. (1974) *Quasi-likelihood functions, Generalized linear Models, and the Gauss-Newton method*. Biometrika, **61**, 439-447.
- Zeger, S.L., Liang, K.Y. (1992) *An overview of methods for the analysis of longitudinal data*. Statistics in medicine, **11**, 1825-1839.
- Zhao, L.P., Prentice, R.L. (1990) *Correlated binary regression using a quadratic exponential model*. Biometrika, **77**, 642-648.